

## logistic regression

$$\hat{y}^{(i)} = \sigma(\omega^T x^{(i)} + b), \text{ where}$$

$$\sigma(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}}$$

$$z^{(i)} = \omega^T x^{(i)} + b$$

$$\left. \begin{matrix} x^{(i)} \\ y^{(i)} \\ z^{(i)} \end{matrix} \right\} \text{ } i\text{th example}$$

Given  $\{ (x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}) \}$ ,  
want  $\hat{y}^{(i)} \approx y^{(i)}$

Loss error function

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

\* if  $y = 1$  :  ~~$L(\hat{y}, y)$~~

$$L(\hat{y}, y) = -(\log \hat{y})$$

In order to get  $L(\hat{y}, y)$  as small as possible,  $\log \hat{y}$  must be large.

But  $\hat{y}$  has limits  $[0, 1]$

so  $\hat{y}$  must be as close to 1 as possible.  $\hat{y} \approx 1$ .

\* if  $y = 0$

$$L(\hat{y}, y) = -\log(1 - \hat{y})$$

~~$$= -\log \hat{y}$$~~

for  $L(\hat{y}, y)$  to be small,  
 $\log(1 - \hat{y})$  must

cost function (On entire training eg)

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

$n \rightarrow$  no. of features

$m \rightarrow$  no. of training examples

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)}) \right]$$

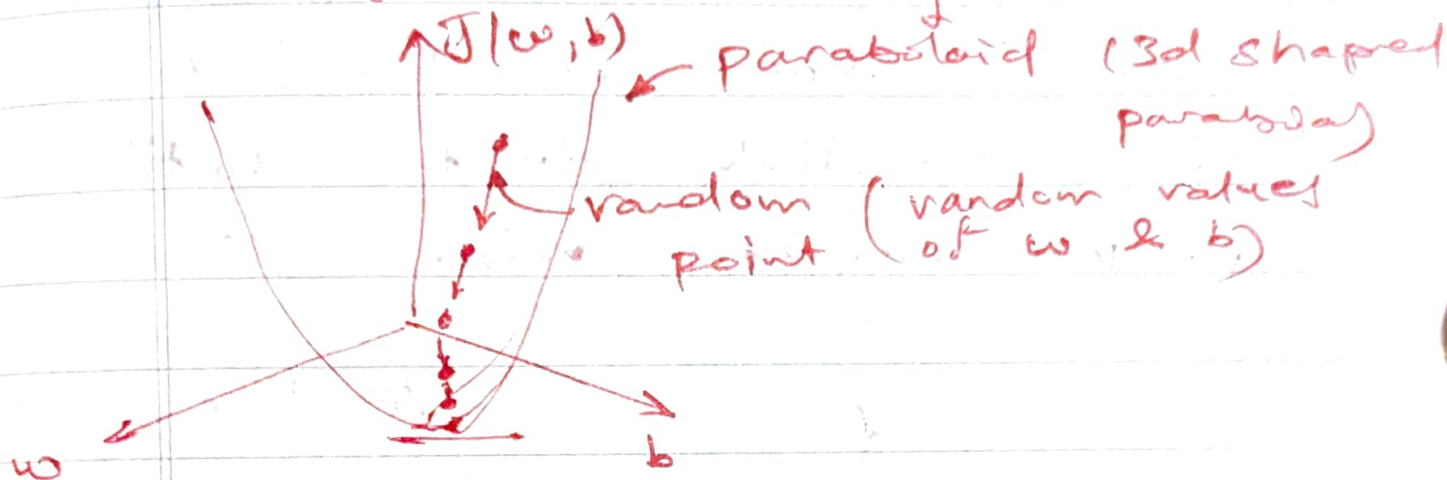
Gradient descent:

what's loss function :- Calculating the loss for one training example

what's cost function :- Calculating the avg loss [i.e. avg of loss function] for one value of  $w$  &  $b$ , (const. value of  $w$  &  $b$ ).

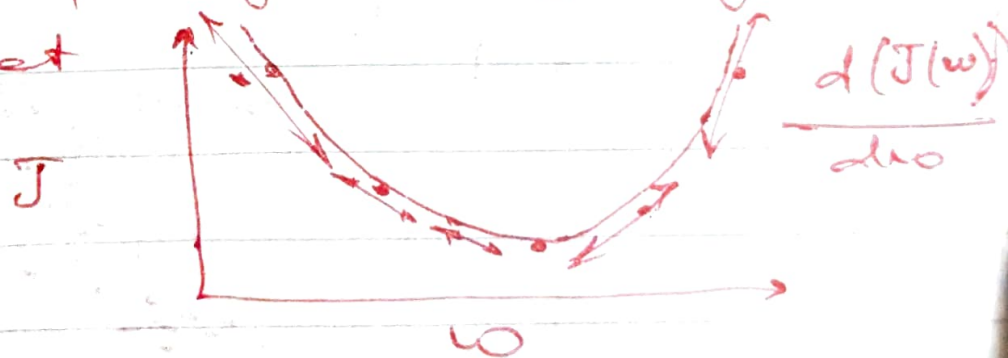
We need to minimize the ~~loss~~ cost function in order to get near perfect output. How can we do this?

Imagine  $w$  and  $b$  are ~~2~~ 1 dimensional. (though  $w$  can be multidimensional, since, training dataset has many features)



We want to reach the minimum of  $J(w, b)$  [to get min. loss].

How can we do this? We can do this by Gradient Descent. If we plot graph neglecting  $b$ , we'll get





We'll calculate the slope at that random point & we'll nudge in the ~~min~~ direction of graph where we'll get min.  $J$ .

This is known as Gradient Descent.

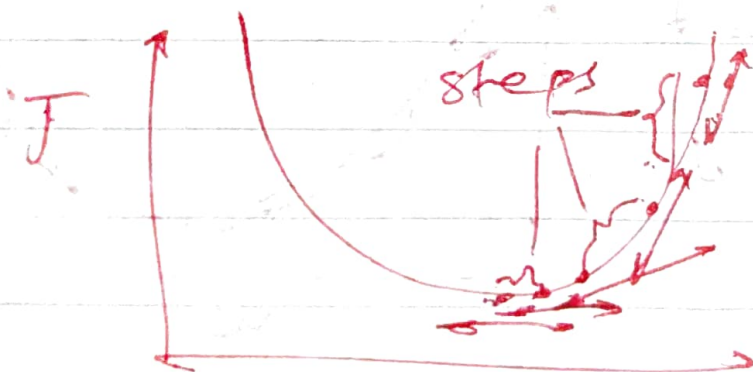
We'll repeatedly do this until we've reached min. of  $J$ .

Repeat {

$$w := w - \alpha \frac{dJ(w)}{dw}$$

learning rate

} denotes  
 $\alpha$  determines the step which we'll have to take to calculate slope at that point.



$w := \dots$       " $:=$ " update sign  
└ update  $w$

No we'll do this for  $b$  as well

Repeat {

$$w := w - \alpha \frac{\delta J(w, b)}{\delta w}$$

$$b := b - \alpha \frac{\delta J(w, b)}{\delta b}$$

}

Since,  $J$  depends on multiple variables, we'll use partial derivative rather than normal derivative

$$a = e^{-z}$$

$$\frac{da}{dz} = \frac{1}{1+e^{-z}}$$

$$= (1+e^{-z})^{-1}$$

$$= \times 1 \cdot (1+e^{-z})^{-2}$$

$$\cdot [0 + \times e^{-z}]$$

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$

~~$$\frac{e^{-z}}{1+2e^{-z}+e^{-2z}}$$~~

~~$$\frac{e^{-z}}{1+e^{-z}}$$~~

$$1-a = 1 - \frac{1}{1+e^{-z}}$$

$$= \frac{1+e^{-z} - 1}{1+e^{-z}}$$

$$= \frac{e^{-z}}{1+e^{-z}}$$

$$da = \frac{\delta L(a, y)}{\delta a}$$

$$= \frac{-y}{a} + \frac{1-y}{1-a} \quad \leftarrow (1)$$

---


$$dz = \frac{\delta L(a, y)}{\delta z}$$

Using chain rule, we can write

$$= \frac{\delta L}{\delta a} \times \frac{\delta a}{\delta z} \quad \leftarrow (2)$$

---


$$\frac{\delta a}{\delta z} = \frac{da}{dz} = \frac{d}{dz} (\sigma(z))$$

$$\sigma(z) = \frac{1}{1+e^{-z}} = (1+e^{-z})^{-1}$$

$$\frac{da}{dz} = -1 (1+e^{-z})^{-2} [0 + (-1)e^{-z}]$$

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$



We know,

$$a = \frac{1}{1+e^{-z}}$$

$$\begin{aligned} 1-a &= 1 - \frac{1}{1+e^{-z}} \\ &= \frac{1+e^{-z} - 1}{1+e^{-z}} \\ &= \frac{e^{-z}}{1+e^{-z}} \end{aligned}$$

$$a(1-a) = \frac{e^{-z}}{(1+e^{-z})^2}$$

$$i \frac{da}{dz} = a(1-a) \quad \text{--- (3)}$$

From (1), (2) & (3)

$$\begin{aligned} dz &= \left( \frac{-y}{a} + \frac{1-y}{1-a} \right) \cdot a(1-a) \\ &= -y(1-a) + (1-y) \cdot a \end{aligned}$$

$$\begin{aligned} &= -y + ay + a - ay \\ &= a - y \end{aligned}$$

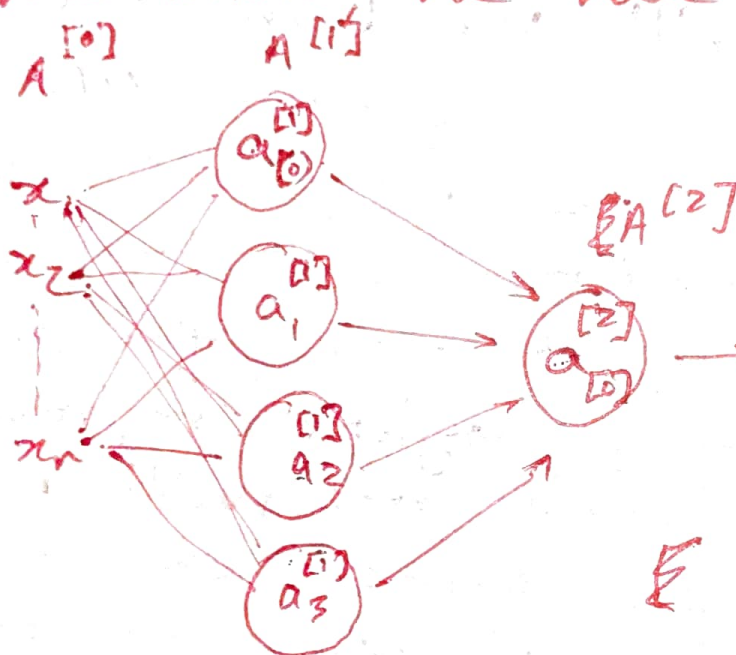
$$\therefore dx = a - y$$

~~(xy)~~

# Vectorizing across multiple examples

no. of test cases / examples =  $m$

We know we have to perform



$[i] \rightarrow$  layer number or activation function layer number

$(j) \rightarrow$  activation function number in that layer

$$\begin{aligned} z &= W^T x + b \\ a &= \sigma(z) \end{aligned}$$

$\rightarrow$  One layer will calc. this & pass "a" to the next layer

Here  $x$  is also input from previous layer,  $\therefore$  input layer.

So, we have to perform

$$z^{[1]} = W^{[0]} A^{[0]} + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[1]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

The  $A^{[2]}$  will result in final output

Since, it's a 2-layer NN.

We can iterate for each example/

training image/test case as it'll cause

a lot of time to train the model.

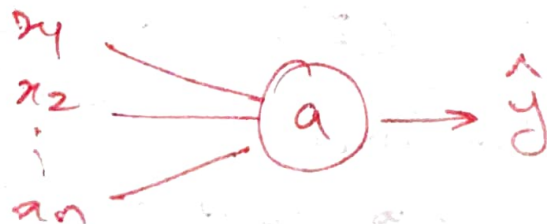
As in previous logistic regression,

we have vectorization, ~~not~~ we'll perform

same here. In previous example though

we had only 1 layer and only 1

activation function & no hidden layer



But here, we have an extra layer  
aka hidden layer & it has more  
than 1 activation function. How  
can we do vectorization of the NN.

Say we have  $m$  training examples & each example has  $n_x$  features.

So  $A^{[0]}$  or  $X$  with should be a matrix of  $(n_x, m)$ . We can do this by stacking  $m$  examples column wise.

In first NN, we have similar matrix  $X$  but here we'll call it as  $A^{[0]}$  as explained before. We'll multiply matrix  $W^{[1]}$  with it, but  $W^{[1]}$  has 4 activation functions. Therefore, dimensions of  $W$  must be  $(4, n_x)$  as

$$W_{4 \times n_x}^{[1]} \times A_{n_x \times m}^{[0]} \quad \text{not no. of}$$

columns ~~and~~ of  $W^{[1]}$  and no. of rows of  $X$  or  $A^{[0]}$  should match. We can achieve these dimensions of  $W^{[1]}$  by stacking up  $n_x$  rows



4 times row with

$$W^{[1]} = \begin{bmatrix} \text{--- } w_{(1)} \text{ ---} \\ \text{--- } w_{(2)} \text{ ---} \\ \text{--- } w_{(3)} \text{ ---} \\ \text{--- } w_{(4)} \text{ ---} \end{bmatrix} \quad 4 \times n_x$$

$$A^{[0]} = \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ x^1 & x^2 & \text{---} & x^m \\ \downarrow & \downarrow & & \downarrow \\ & & & \downarrow \\ & & & n_x \times n_h \end{bmatrix}$$

training examples

and we can add  $b$  to it.  
We don't need to specify its dimensions as python will perform broadcasting for it.

$$A^{[1]} = \sigma(z^{[1]}) = W^{[1]} \cdot A^{[0]} + b^{[1]}$$

After calculating  $z^{[1]}$ , we can calculate  $A^{[1]}$ . Its dimensions would ~~with~~

be  $4 \times m$  (as  $W^{[1]}_{4 \times m} \times A^{[0]}_{m \times m}$ )  
~~So~~

Now, we need to calculate  $A^{[2]}$   
 $A^{[2]} = \sigma(z^{[2]}) = W^{[2]} A^{[1]} + b^{[2]}$

What should be the dimensions for  $W^{[2]}$ .

$A^{[1]}$  has  $4 \times m$  dimensions.

So  $W^{[2]}$  must have 4 columns.

What about rows?

As  $A^{[2]}$  is directly responsible for the output i.e.  $\hat{y}$ , it must be an unidimensional vector.  $A^{[2]}_{1 \times 1}$

~~$A^{[2]}$~~  In order to get that,

$W^{[2]}$  must have only 1 row.

$\therefore$  Dimensions of  $W^{[2]} = (1 \times 4)$

$b^{[2]}$  will be a real number.

~~$A^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$~~

$\therefore A^{[2]} = \sigma(z^{[2]}) = W^{[2]} A^{[1]} + b^{[2]}$

# Neural Network representation learning

$$z^{[1]} = W^{[1]} x + b^{[1]}$$

$4 \times 1$                        $(4, 3)$      $(3, 1)$      $(4, 1)$

$\nearrow$  or  $a^{[0]}$

$$a^{[1]} = \sigma(z^{[1]})$$

$(4, 1)$                        $(4, 1)$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$(1, 1)$                        $(1, 4)$      $(4, 1)$      $+ (1, 1)$

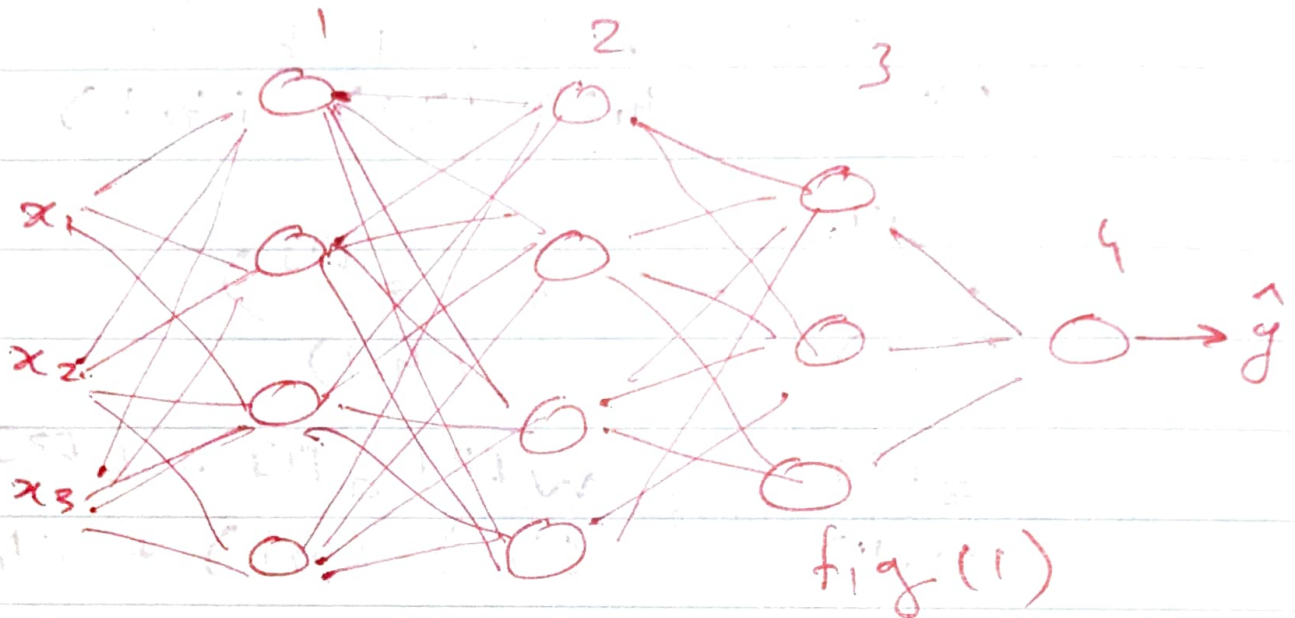
$$a^{[2]} = \sigma(z^{[2]})$$

$(1, 1)$                        $(1, 1)$

↑ This is for 1 training example

Note: A → for m examples  
a → for 1 example

Forward propagation for multilayer NN (aka deep network)



forward propagation:

$$z^{[1]} = W^{[1]} A^{[0]} + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$z^{[4]} = W^{[4]} A^{[3]} + b^{[4]}$$

$$A^{[4]} = \hat{y} = \sigma(z^{[4]})$$



general representation

$$z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

vectorize it over  $l$  layers.

for  $l = 1 ; l \leq 5 ; l++ :$

$$z^{[l]} = w^{[l]} A^{[l-1]}$$

$$z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

Dimensions of  $z, w, b, dw, db$  throughout NN for diff. layers.

for fig(1), we have 3 input features

$$z^{[1]} = w^{[1]} \cdot x + b^{[1]}$$

$x$  has dimensions  $(3, 1)$  — for 1 training example  
 $z^{[1]}$  has dim.  $= (4, 1)$  — since it has 4 nodes

$$\begin{array}{ccccc} z^{[1]} & = & w^{[1]} & \cdot & x & + & b^{[1]} \\ (4, 1) & & (4, 3) & & (3, 1) & & (4, 1) \end{array}$$

$w^{[1]}$  must have dim.  $(4, 3)$



We can say,  $w^{[1]}$  has dimensions  
 $(n^{[1]}, n^{[0]})$   $n^{[1]} \rightarrow$  no. of nodes in  
layer 1.

more generally, we can say,

$$w^{[L]} = (n^{[L]}, n^{[L-1]})$$

$$z^{[L]} = (n^{[L]}, 1)$$

$$b^{[L]} = (n^{[L]}, 1)$$

$$a^{[L]} = g^{[L]}(z^{[L]})$$

$$\therefore \text{dim. of } A^{[L]} = (n^{[L]}, 1)$$

dim. of  $d_w$  &  $d_b$  ~~are~~ should be the  
same as  $w$  &  $b$  resp.

} for 1 training  
example